

# Practical Application III: Comparing Classifiers

**Overview:** In this practical application, your goal is to compare the performance of the classifiers we encountered in this section, namely K Nearest Neighbor, Logistic Regression, Decision Trees, and Support Vector Machines. We will utilize a dataset related to marketing bank products over the telephone.

## Getting Started

Our dataset comes from the UCI Machine Learning repository [link](#). The data is from a Portuguese banking institution and is a collection of the results of multiple marketing campaigns. We will make use of the article accompanying the dataset [here](#) for more information on the data and features.

## Problem 1: Understanding the Data

To gain a better understanding of the data, please read the information provided in the UCI link above, and examine the **Materials and Methods** section of the paper. How many marketing campaigns does this data represent?

In [593]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from ipynb_utils import show_html_output
import seaborn as sns
import plotly.express as px
from sklearn import preprocessing
from sklearn.datasets import load_digits
from sklearn.preprocessing import LabelEncoder
from sklearn.model_selection import train_test_split, GridSearchCV
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.dummy import DummyClassifier
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import RandomForestClassifier
from sklearn.neighbors import KNeighborsClassifier
from time import time
from sklearn.metrics import confusion_matrix
from sklearn.metrics import ConfusionMatrixDisplay
from sklearn.svm import SVC
```

In [594]:

```
df = pd.read_csv('data/bank-additional-full.csv', sep = ';')
```

In [595]:

```
df.head()
```

Out[595]:

	age	job	marital	education	default	housing	loan	contact	month	day_of_week	...	campaign	pdays	previous	poutcome
0	56	housemaid	married	basic.4y	no	no	no	telephone	may	mon	...	1	999	0	noneexistent
1	57	services	married	high.school	unknown	no	no	telephone	may	mon	...	1	999	0	noneexistent
2	37	services	married	high.school	no	yes	no	telephone	may	mon	...	1	999	0	noneexistent
3	40	admin.	married	basic.6y	no	no	no	telephone	may	mon	...	1	999	0	noneexistent
4	56	services	married	high.school	no	no	yes	telephone	may	mon	...	1	999	0	noneexistent

5 rows × 21 columns

## Problem 3: Understanding the Features

Examine the data description below, and determine if any of the features are missing values or need to be coerced to a different data type.

**Input variables:**

- 1 – bank client data:
- 2 – age (numeric)
- 3 – job : type of job (categorical: 'admin.', 'blue-collar', 'entrepreneur', 'housemaid', 'management', 'retired', 'self-employed', 'services', 'student', 'technician', 'unemployed', 'unknown')
- 4 – marital : marital status (categorical: 'divorced', 'married', 'single', 'unknown'; note: 'divorced' means divorced or widowed)
- 5 – education (categorical: 'basic.4y', 'basic.6y', 'basic.9y', 'high.school', 'illiterate', 'professional.course', 'university.degree')
- 6 – default: has credit in default? (categorical: 'no', 'yes', 'unknown')
- 7 – housing: has housing loan? (categorical: 'no', 'yes', 'unknown')
- 8 – loan: has personal loan? (categorical: 'no', 'yes', 'unknown')
- 9 – contact with the last contact of the current campaign:
- 10 – contact: contact communication type (categorical: 'cellular', 'telephone')
- 11 – month: last contact month of year (categorical: 'jan', 'feb', 'mar', ..., 'nov', 'dec')
- 12 – day\_of\_week: last contact day of the week (categorical: 'mon', 'tue', 'wed', 'thu', 'fri')
- 13 – duration: last contact duration, in seconds (numeric). Important note: this attribute highly affects the output target (e.g., if duration=0 then y=0). Yet, the duration is not known before a call is performed. Also, after the end of the call y is obviously known. Thus, this input should only be included for benchmark purposes and should be discarded if the intention is to have a realistic predictive model.
- 14 – campaign: number of contacts performed during this campaign and for this client (numeric, includes last contact)
- 15 – pdays: number of days that passed by after the client was last contacted from a previous campaign (numeric; 999 means client was not previously contacted)
- 16 – previous: number of contacts performed before this campaign and for this client (numeric)
- 17 – poutcome: outcome of the previous marketing campaign (categorical: 'failure', 'nonexistent', 'success')
- 18 – cons.price.idx: consumer price index – quarterly indicator (numeric)
- 19 – cons.conf.idx: consumer confidence index – monthly indicator (numeric)
- 20 – euribor3m: euribor 3 month rate – daily indicator (numeric)
- 21 – nr.employed: number of employees (numeric)
- 22 – y – has the client subscribed a term deposit? (binary: 'yes', 'no')

In [596]:

```
df.isnull().sum()
```

Out[596]:

age	0
job	0
marital	0
education	0
default	0
housing	0
loan	0
contact	0
month	0
day_of_week	0
duration	0
campaign	0
pdays	0
previous	0
poutcome	0
emp.var.rate	0
cons.price.idx	0
cons.conf.idx	0
euribor3m	0
y	0
dtype: int64	

## Problem 4: Understanding the Task

After examining the description and data, your goal now is to clearly state the Business Objective of the task. State the objective below.

In [598]:

```
df.info()
```

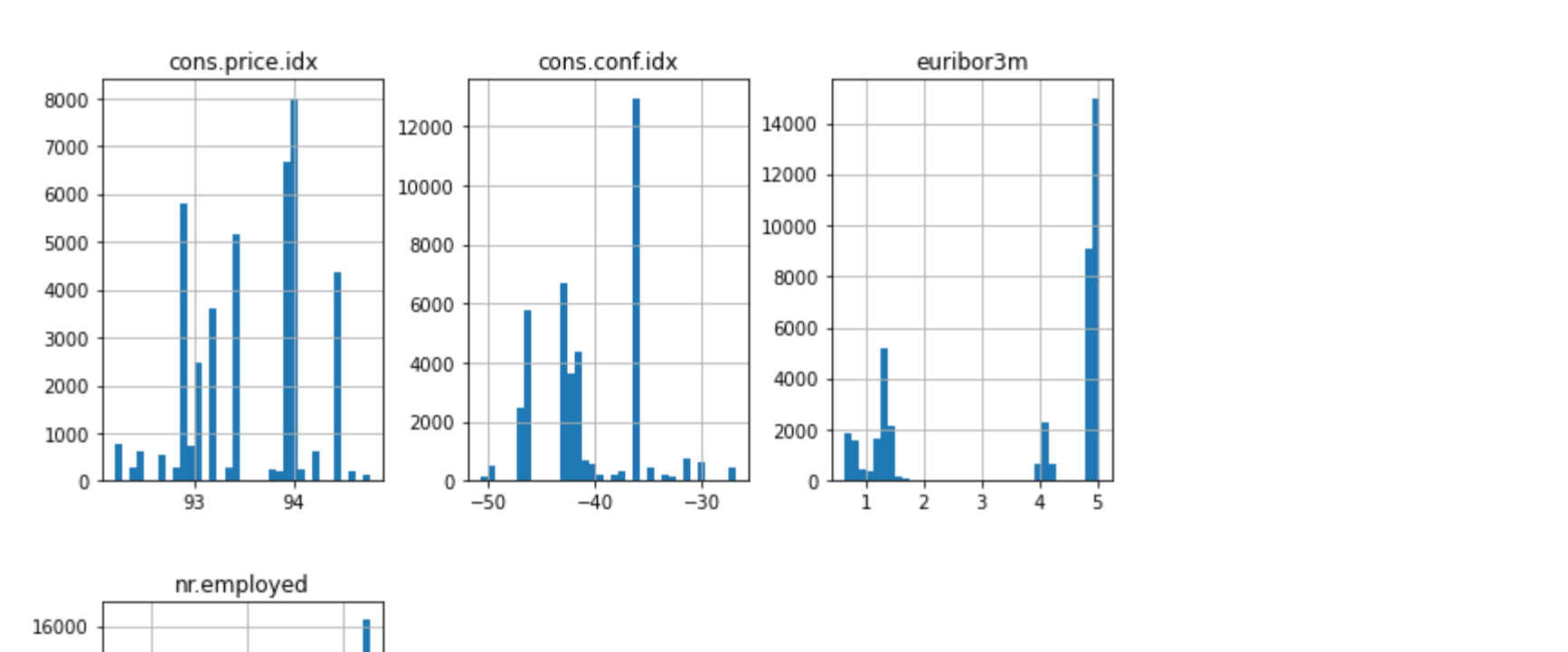
<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 41188 entries, 0 to 41187  
Data columns (total 21 columns):  
# Column Non-Null Count Dtype  
---  
0 age 41188 non-null int64  
1 job 41188 non-null object  
2 marital 41188 non-null object  
3 education 41188 non-null object  
4 default 41188 non-null object  
5 housing 41188 non-null object  
6 loan 41188 non-null object  
7 contact 41188 non-null object  
8 month 41188 non-null object  
9 day\_of\_week 41188 non-null object  
10 duration 41188 non-null int64  
11 campaign 41188 non-null int64  
12 pdays 41188 non-null int64  
13 previous 41188 non-null int64  
14 poutcome 41188 non-null object  
15 emp.var.rate 41188 non-null float64  
16 cons.price.idx 41188 non-null float64  
17 cons.conf.idx 41188 non-null float64  
18 euribor3m 41188 non-null float64  
19 nr.employed 41188 non-null float64  
20 y 41188 non-null object  
dtypes: float64(5), int64(5), object(11)  
memory usage: 6.6+ MB

In [599]:

```
df.hist(bins=35, figsize=(10,20))
```

Out[599]:

array([[<matplotlib.axes.\_subplots.AxesSubplot object at 0x156905040>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x156932310>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x15693a2c0>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x1564f280>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x156c7a00>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x156c86100>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x156c861f0>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x156bd9f00>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x1567d0800>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x15f8030a0>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x15d497000>, <matplotlib.axes.\_subplots.AxesSubplot object at 0x15d497070>],  
dtype=object)



In [600]:

```
# create Correlation Matrix to understand the correlation between numeric columns  
df.corr()
```

<ipython-input-600-0789589e9e74>:12: FutureWarning: The default value of numeric\_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric\_only to silence this warning.

Out[600]:

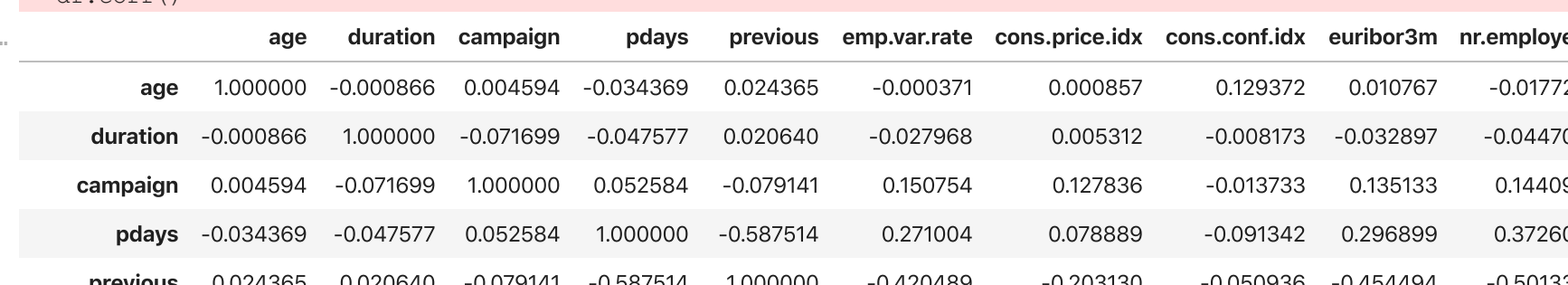
	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons.conf.idx	euribor3m	nr.employed
age	1.000000	-0.000866	1.000000	-0.047699	-0.034369	0.024365	-0.000371	0.129372	0.010767	-0.017725
duration	0.0004594	1.000000	-0.000000	-0.047677	0.020640	-0.027968	0.005312	-0.008173	-0.032897	-0.044703
campaign	0.0004360	-0.007699	1.000000	-0.000000	-0.057514	0.027004	0.078889	-0.091342	0.296989	-0.372605
pdays	0.024365	-0.007677	-0.007941	1.000000	-0.587514	0.000000	-0.420489	-0.203130	-0.050936	-0.454494
previous	0.000857	0.005312	0.078889	-0.000000	1.000000	0.775334	0.196041	0.972245	0.906970	0.906970
emp.var.rate	-0.000371	-0.027968	0.027004	-0.420489	0.000000	1.000000	0.058986	0.688230	0.522034	0.100513
cons.price.idx	0.129372	-0.008173	-0.091342	-0.203130	0.775334	0.196041	1.000000	0.277866	1.000000	0.945154
cons.conf.idx	0.010767	-0.032897	0.296989	-0.050936	0.906970	0.972245	0.688230	1.000000	1.000000	0.945154
euribor3m	-0.017725	-0.044703	0.372605	-0.501333	0.906970	0.522034	0.100513	0.945154	1.000000	1.000000
nr.employed	-0.017725	-0.044703	0.372605	-0.501333	0.906970	0.522034	0.100513	0.945154	1.000000	1.000000

In [601]:

```
fig, ax = plt.subplots(figsize=(20,10))  
sns.heatmap(df.corr(), annot=True)
```

<ipython-input-601-21f4ec2ce54>:12: FutureWarning: The default value of numeric\_only in DataFrame.corr is deprecated. In a future version, it will default to False. Select only valid columns or specify the value of numeric\_only to silence this warning.

Out[601]:



In [602]:

```
df.info()
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 41188 entries, 0 to 41187  
Data columns (total 21 columns):  
# Column Non-Null Count Dtype  
---  
0 age 41188 non-null int64  
1 job 41188 non-null object  
2 marital 41188 non-null object  
3 education 41188 non-null object  
4 default 41188 non-null object  
5 housing 41188 non-null object  
6 loan 41188 non-null object  
7 contact 41188 non-null object  
8 month 41188 non-null object  
9 day\_of\_week 41188 non-null object  
10 duration 41188 non-null int64  
11 campaign 41188 non-null int64  
12 pdays 41188 non-null int64  
13 previous 41188 non-null int64  
14 poutcome 41188 non-null object  
15 emp.var.rate 41188 non-null float64  
16 cons.price.idx 41188 non-null float64  
17 cons.conf.idx 41188 non-null float64  
18 euribor3m 41188 non-null float64  
19 nr.employed 41188 non-null float64  
20 y 41188 non-null object  
dtypes: float64(5), int64(5), object(11)  
memory usage: 6.6+ MB

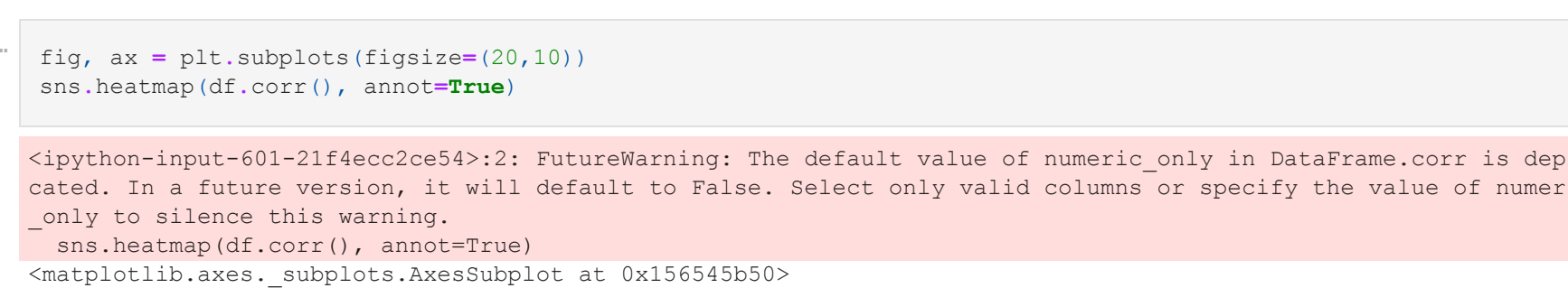
In [603]:

```
num_col = ['age', 'duration', 'campaign', 'pdays', 'previous', 'emp.var.rate', 'cons.price.idx', 'cons.conf.idx']  
cat_col = ['job', 'marital', 'education', 'default', 'housing', 'loan', 'contact', 'month', 'day_of_week', 'poutcome']
```

In [604]:

```
df["age"].hist(bins=10)
```

Out[604]:



In [605]:

```
print(df["age"].max())  
print(df["age"].min())
```

Out[605]:

In [606]:

```
bank_df = df.copy()
```

In [607]:

```
bank_df.info()
```

<class 'pandas.core.frame.DataFrame'>  
RangeIndex: 41188 entries, 0 to 41187  
Data columns (total 21 columns):  
# Column Non-Null Count Dtype  
---  
0 age 41188 non-null int64  
1 job 41188 non-null object  
2 marital 41188 non-null object  
3 education 41188 non-null object  
4 default 41188 non-null object  
5 housing 41188 non-null object  
6 loan 41188 non-null object  
7 contact 41188 non-null object  
8 month 41188 non-null object  
9 day\_of\_week 41188 non-null object  
10 duration 41188 non-null int64  
11 campaign 41188 non-null int64  
12 pdays 41188 non-null int64  
13 previous 41188 non-null int64  
14 poutcome 41188 non-null object  
15 emp.var.rate 41188 non-null float64  
16 cons.price.idx 41188 non-null float64  
17 cons.conf.idx 41188 non-null float64  
18 euribor3m 41188 non-null float64  
19 nr.employed 41188 non-null float64  
20 y 41188 non-null object  
dtypes: float64(5), int64(5), object(11)  
memory usage: 6.6+ MB

In [608]:

```
bank_df.describe()
```

Out[608]:

	age	duration	campaign	pdays	previous	emp.var.rate	cons.price.idx	cons.conf.idx	euribor3m
count	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000	41188.000000
mean	40.22406	258.285010	2.567593	986.475454	0.172963	0.061886	93.575664	-40.502600	3.621291
std	10.421225	259.729249	2.777004	186.910907	0.494901	1.570960	0.578840	4.628198	1.734447
min	17.000000	0.000000	1.000000	0.000000	0.000000	-3.400000	92.201000	-52.700000	0.634000
25%	32.000000	102.000000	1.000000	999.000000	0.000000	-1.800000	93.075000	-42.700000	1.344000
50%	38.000000	180.000000	2.000000	999.000000	0.000000	1.000000	93.749000	-41.800000	4.857000
75%	47.000000	319.000000	3.000000	999.000000	0.000000	1.400000	93.994000	-36.400000	4.961000
max	98.000000	4918.000000	56.000000	999.000000	7.000000	1.400000	94.767000	-26.900000	5.045000

In [609]:

```
# Removing the outliers  
first_quartile = df["age"].quantile(.25)  
third_quartile = df["age"].quantile(.75)  
iqr = third_quartile - first_quartile  
lower_quartile = first_quartile - 1.5 * iqr  
upper_quartile = third_quartile + 1.5 * iqr  
bank_df = df.loc[(df["age"] > lower_quartile) & (df["age"] < upper_quartile)]
```

In [610]:

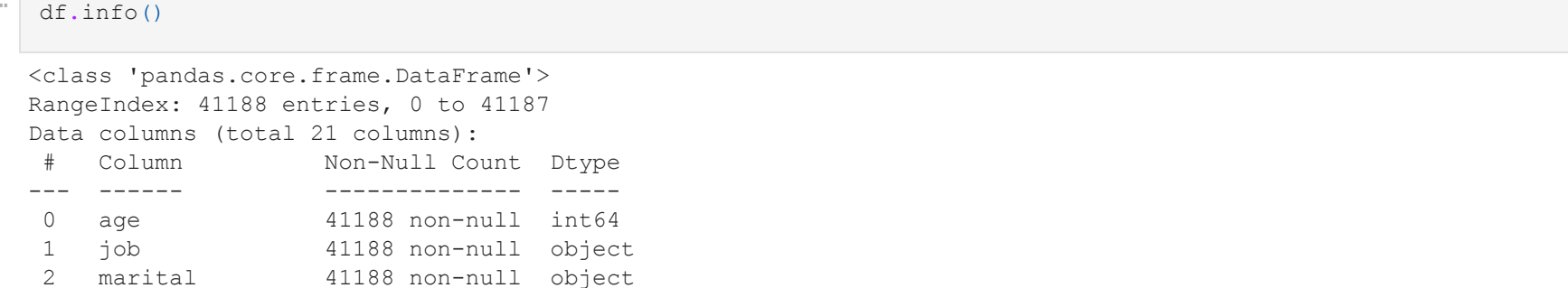
```
print(bank_df["age"].max())  
print(bank_df["age"].min())
```

Out[610]:

In [611]:

```
df["duration"].hist(bins=10)
```

Out[611]:



In [612]:

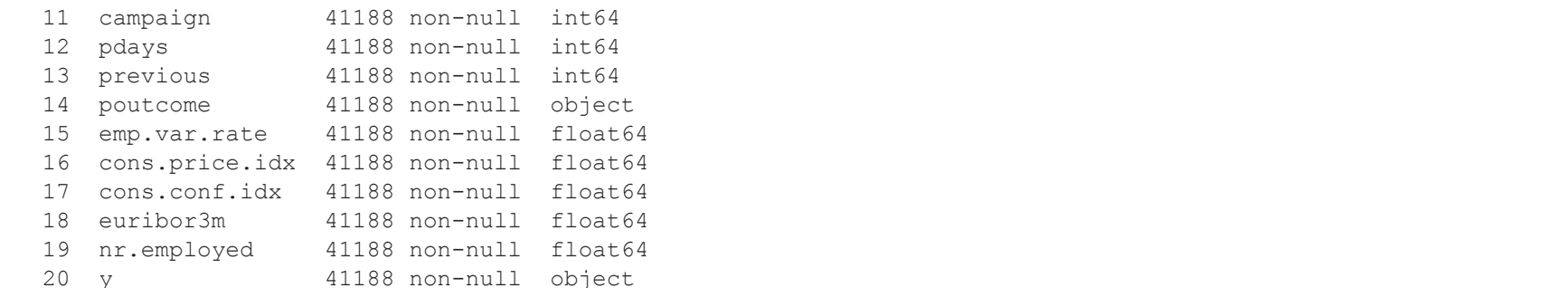
```
# Removing outliers from duration  
first_quartile = bank_df["duration"].quantile(.25)  
third_quartile = bank_df["duration"].quantile(.75)  
iqr = third_quartile - first_quartile  
lower_quartile = first_quartile - 1.5 * iqr  
upper_quartile = third_quartile + 1.5 * iqr  
bank_df = df.loc[(df["duration"] > lower_quartile) & (df["duration"] < upper_quartile)]
```

Out[612]:

In [613]:

```
bank_df["duration"].hist(bins=10)
```

Out[613]:



In [614]:

```
print(bank_df["campaign"].max())  
print(bank_df["campaign"].min())
```

Out[614]:

In [615]:

```
# Remove outliers from the numerical column - campaign using IQR  
first_quartile = bank_df["campaign"].quantile(.25)  
third_quartile = bank_df["campaign"].quantile(.75)  
iqr = third_quartile - first_quartile  
lower_quartile = first_quartile - 1.5 * iqr  
upper_quartile = third_quartile + 1.5 * iqr  
bank_df = df.loc[(df["campaign"] > lower_quartile) & (df["campaign"] < upper_quartile)]  
bank_df["campaign"].hist(bins=10)
```

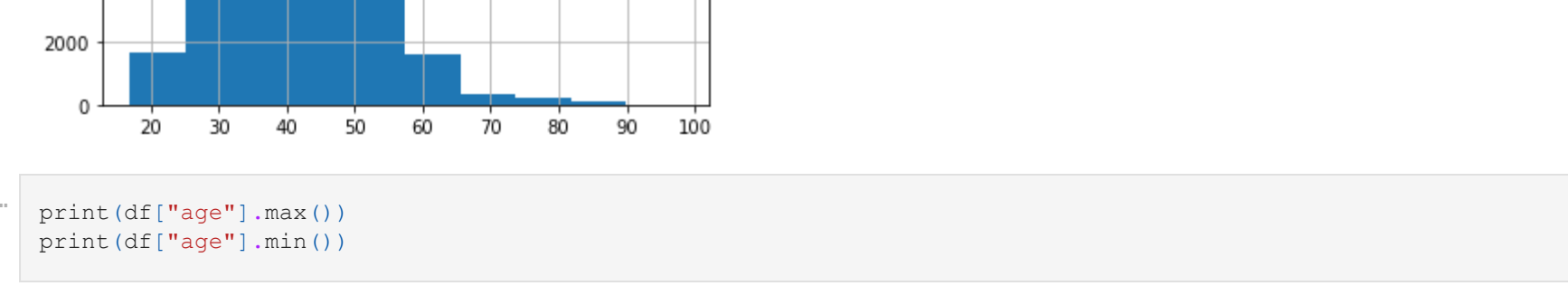
Out[615]:



In [616]:

```
bank_df.hist(bins=30, figsize=(12, 10))
```

Out[616]:

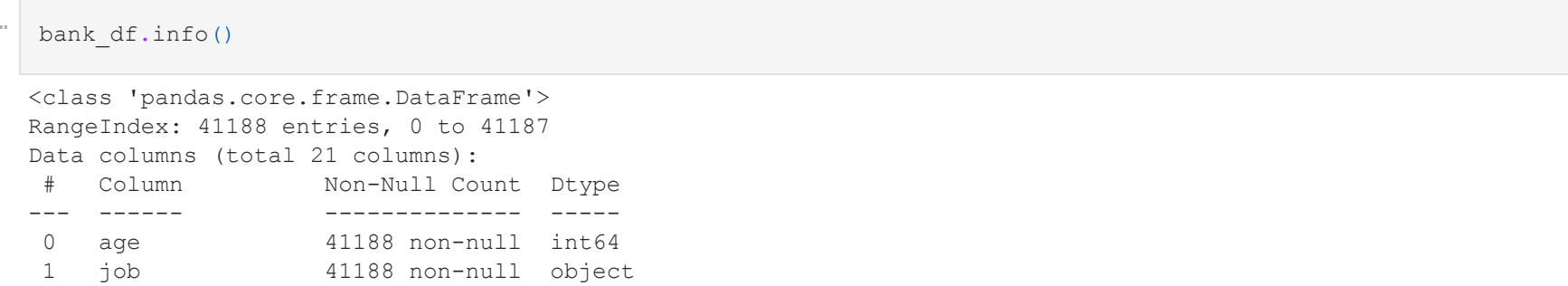


In [617]:

```
# Campaign Acceptance rate based of Number of Contact performed  
fig, ax = plt.subplots(figsize=(10,5))  
sns.countplot(data=bank_df, x="campaign", hue="y")  
plt.setp(plt.xticks()[1:], rotation=45)  
plt.title("Campaign Acceptance Rate based on Number of Contacts Performed")
```

Out[617]:

Text(0.5, 1.0, 'Campaign Acceptance Rate based on Number of Contacts Performed')

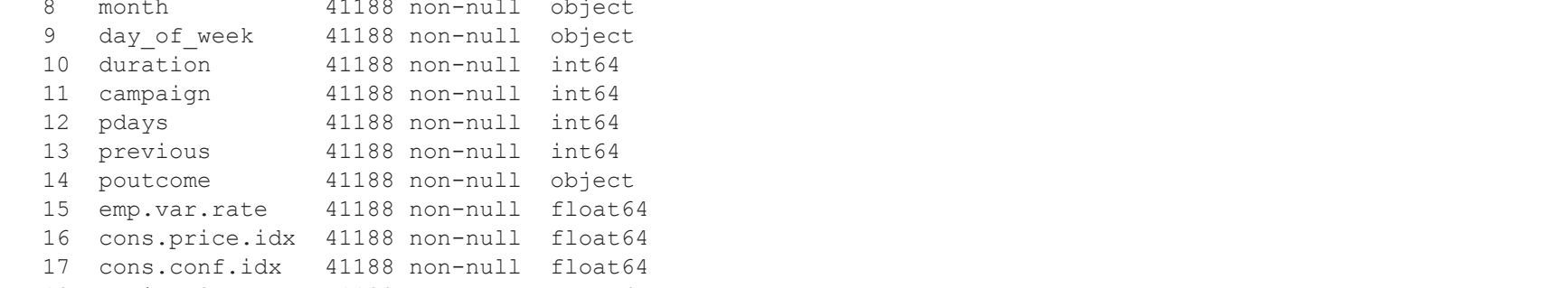


In [618]:

```
# Campaign Acceptance rate based of Jobs  
fig, ax = plt.subplots(figsize=(10,5))  
plt.setp(plt.xticks()[1:], rotation=45)  
sns.countplot(data=bank_df, x="job", hue="y")  
plt.title("Campaign Acceptance Rate based on Jobs Performed")
```

Out[618]:

Text(0.5, 1.0, 'Campaign Acceptance Rate based on Jobs Performed')

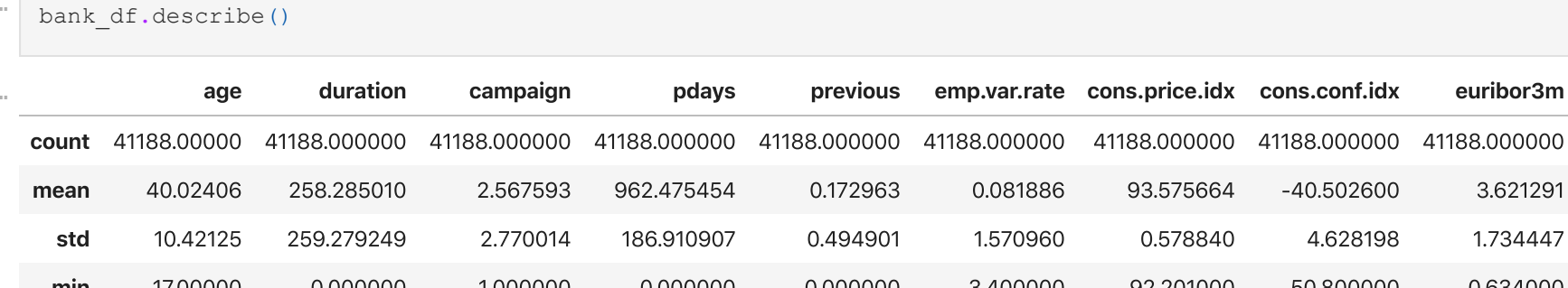


In [620]:

```
# Campaign Acceptance rate based of Marital Status  
fig, ax = plt.subplots(figsize=(10,5))  
plt.setp(plt.xticks()[1:], rotation=45)  
sns.countplot(data=bank_df, x="marital", hue="y")  
plt.title("Campaign Acceptance Rate based on Marital Status")
```

Out[620]:

Text(0.5, 1.0, 'Campaign Acceptance Rate based on Marital Status')

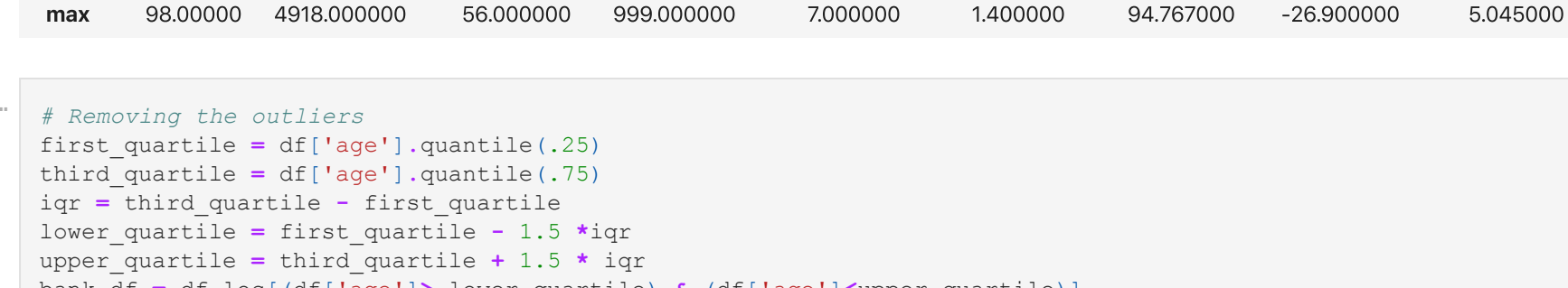


In [621]:

```
# Campaign Acceptance rate based of Education  
fig, ax = plt.subplots(figsize=(10,5))  
plt.setp(plt.xticks()[1:], rotation=45)  
sns.countplot(data=bank_df, x="education", hue="y")  
plt.title("Campaign Acceptance Rate based on Education")
```

Out[621]:

Text(0.5, 1.0, 'Campaign Acceptance Rate based on Education')



In [622]:

```
# Campaign Acceptance rate based of Education  
fig, ax = plt.subplots(figsize=(10,5))  
plt.setp(plt.xticks()[1:], rotation=45)  
sns.countplot(data=bank_df, x="loan", hue="y")  
plt.title("Campaign Acceptance Rate based on loan")
```

Out[622]:

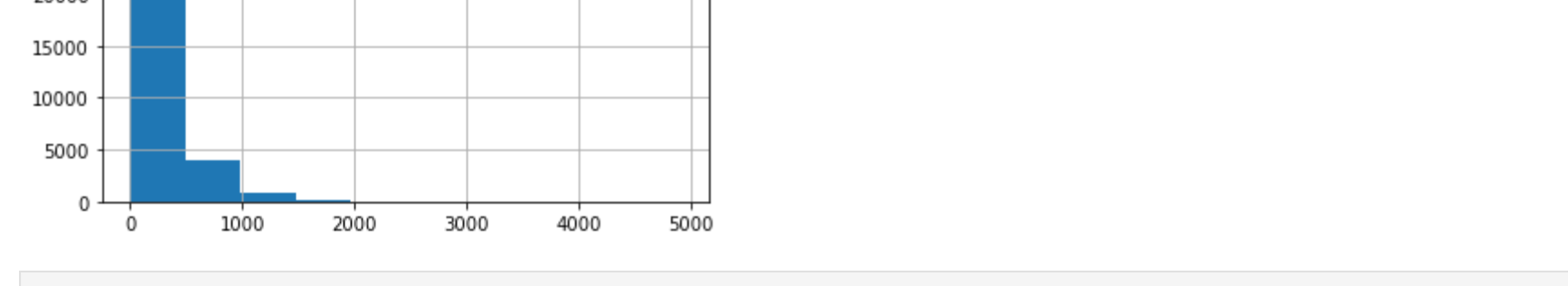
Text(0.5, 1.0, 'Campaign Acceptance Rate based on loan')



In [623]:

```
# Campaign Acceptance rate based of Education  
fig, ax = plt.subplots(figsize=(10,5))  
plt.setp(plt.xticks()[1:], rotation=45)  
sns.countplot(data=bank_df, x="poutcome", hue="y")  
plt.title("Campaign Acceptance Rate based on outcome")
```

Out[623]:



## Problem 5: Engineering Features

Now that you understand your business objective, we will build a basic model to get started. Before we can do this, we must work to encode the data. Using just the bank information features (columns 1 - 7), prepare the features and target column for modeling with appropriate encoding and transformations.

In [624]:

```
from sklearn.preprocessing import LabelEncoder  
encoder = preprocessing.LabelEncoder()  
encoder_df = bank_df.copy()  
cat_col = encoder.fit_transform(encoder_df[cat_col])  
impute_ordinal = encoder.fit_transform(data)  
data.loc[data.notnull()] = np.squeeze(impute_ordinal)  
return data
```

Out[624]:

0.0 [1] [00:00:07, 717/s]<ipython-input-624-a6fc6193fbfe>:6: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

100% [1] [00:00:00, 113.71/s]<ipython-input-624-a6fc6193fbfe>:6: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame



```
In [625.]: count      37803.000000    37803.000000    37803.000000    37803.000000    37803.000000    37803.000000    37803.000000    37803.000000    37803.000000
mean      40.007142      263.432717      194.1751      960.176653      0.181785      0.028646      93.560179      -40.485398      3.5691
std       10.458787      259.251434      1125.0771      192.472659      0.506630      1.553386      0.581213      4.693541      1.7483
min       17.000000      0.000000      1.000000      0.000000      0.000000      -3.400000      92.201000      -50.800000      0.6340
50%       32.000000      107.000000      1.000000      999.000000      0.000000      -1.800000      93.075000      -42.700000      1.3340
75%       38.000000      185.000000      2.000000      999.000000      0.000000      1.000000      93.444000      -41.800000      4.8570
max       47.000000      324.000000      3.000000      999.000000      0.000000      1.400000      93.994000      -36.400000      9.9610
mean      98.800000      4918.000000      5.000000      999.000000      7.000000      1.400000      94.767000      -26.900000      5.0450

In [626.]: encoder_df = encoder_df.drop('duration', axis=1)

In [627.]: encoder_df.reset_index()

Out[627.]:   index  age  job  marital  education  default  housing  loan  contact  month  ...  campaign  pdays  previous  outcome  emp.var.rate
0      0    56   3      1      0      0      0      0      0      1      6  ...      1      999      0      1
1      1    57   7      1      3      1      0      0      0      1      6  ...      1      999      0      1
2      2    37   7      1      3      0      2      0      1      6  ...      1      999      0      1
3      3    40   0      1      1      0      0      0      1      6  ...      1      999      0      1
4      4    56   7      1      3      0      2      1      6  ...      1      999      0      1
...    ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
37798  41183  73   5      1      5      0      2      0      0      7  ...      1      999      0      1
37799  41184  46   1      1      5      0      0      0      0      7  ...      1      999      0      1
37800  41185  56   5      1      6      0      2      0      0      7  ...      2      999      0      1
37801  41186  44   9      1      5      0      0      0      0      7  ...      1      999      0      1
37802  41187  74   5      1      5      0      2      0      0      7  ...      3      999      1      0

37803 rows x 21 columns

In [628.]: #standardizing the dataset
scaler = preprocessing.StandardScaler()
std_df = scaler.fit_transform(encoder_df.drop('job', 'marital', 'education', 'default', 'housing', 'loan', 'contact', 'month', 'previous', 'outcome', 'emp.var.rate', 'n.employed', 'n.employed_rate'))

In [629.]: scaled_df = pd.DataFrame(std_df, columns=['age', 'pdays', 'cons.price.idx', 'cons.conf.idx', 'euribor3m', 'nr.employed'])

In [630.]: scaled_df['job'] = encoder_df['job'].values
scaled_df['marital'] = encoder_df['marital'].values
scaled_df['education'] = encoder_df['education'].values
scaled_df['default'] = encoder_df['default'].values
scaled_df['housing'] = encoder_df['housing'].values
scaled_df['loan'] = encoder_df['loan'].values
scaled_df['contact'] = encoder_df['contact'].values
scaled_df['month'] = encoder_df['month'].values
scaled_df['day_of_week'] = encoder_df['day_of_week'].values
scaled_df['campaign'] = encoder_df['campaign'].values
scaled_df['previous'] = encoder_df['previous'].values
scaled_df['outcome'] = encoder_df['outcome'].values
scaled_df['emp.var.rate'] = encoder_df['emp.var.rate'].values
scaled_df['y'] = encoder_df['y'].values

In [631.]: scaled_df.head()

Out[631.]:   age  pdays  cons.price.idx  cons.conf.idx  euribor3m  nr.employed  job  marital  education  default  housing  loan  contact  month  previous  outcome
0  1.529151  0.201712  0.746417  0.870441  0.736639  0.360955  3  1  0  0  0  0  0  1
1  1.624766  0.201712  0.746417  0.870441  0.736639  0.360955  7  1  3  0  1  0  0  1
2 -0.287527  0.201712  0.746417  0.870441  0.736639  0.360955  7  1  3  0  2  0  1
3 -0.005853  0.201712  0.746417  0.870441  0.736639  0.360955  0  1  1  0  0  0  1
4  1.529151  0.201712  0.746417  0.870441  0.736639  0.360955  7  1  3  0  0  2  1

In [632.]: scaled_df.describe()

Out[632.]:   age  pdays  cons.price.idx  cons.conf.idx  euribor3m  nr.employed  campaign  previous  emp.var.rate
count  3.780300e+04  3.780300e+04  3.780300e+04  3.780300e+04  3.780300e+04  3.780300e+04  37803.000000  37803.000000  37803.000000
mean  3.247939e-16  -2.17639e-17  -1.814033e-14  0.000000  2.405680e-17  4.811759e-16  1.941751  0.181785  0.028646
std  1.000013e+00  1.000013e+00  1.000013e+00  1.000013  1.000013e+00  1.000013e+00  1.125077  0.506630  1.553386
min  -2.199820e+00  -4.988721e+00  -2.338552e+00  -2.197645  -1.678834e+00  -2.752819e+00  1.000000  0.000000  -3.400000
50%  -7.65601e-01  2.017117e-01  -8.347802e-01  -0.471847  -1.278447e+00  -5.974258e-01  1.000000  0.000000  1.000000
75%  6.686196e-01  2.017117e-01  7.464167e-01  0.870441  7.961253e-01  8.689632e-01  3.000000  0.000000  1.400000
max  5.544967e+00  2.017117e-01  2.076411e+00  2.894526  8.441717e-01  8.689632e-01  5.000000  7.000000  9.961000

Problem 6: Train/Test Split
With your data prepared, split it into a train and test set.

In [633.]: X = scaled_df.drop('y', axis=1)
y = scaled_df['y'].astype('int')
print(X.shape)
print(y.shape)

(37803, 19)
(37803,)

In [634.]: X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.33, random_state=42)

In [635.]: print(X_train.shape)
print(y_train.shape)
print(X_test.shape)
print(y_test.shape)

(25328, 19)
(25328,)
(12475, 19)
(12475,)

In [ ]:

Problem 7: A Baseline Model
Before we build our first model, we want to establish a baseline. What is the baseline performance that our classifier should aim to beat?

In [636.]: models = ['Baseline']
train_time = []
train_accuracy = []
test_accuracy = []
accuracy_score = []
AUC_score = []

In [637.]: # Baseline Model
dummy_clf = DummyClassifier(random_state=42)
start_time = time()
dummy_clf.fit(X_train, y_train)
train_time.append(time() - start_time)
# Training
y_preds = dummy_clf.predict(X_train)
train_accuracy_score = dummy_clf.score(X_train, y_train)
train_accuracy.append(train_accuracy_score)
# Test score
y_preds_test = dummy_clf.predict(X_test)
test_score = dummy_clf.score(X_test, y_test)
test_accuracy.append(test_accuracy_score)
accuracy_score.append((train_accuracy_score + test_accuracy_score) / 2)
AUC_score.append('N/A')

Problem 8: A Simple Model
Use Logistic Regression to build a basic model on your data.

In [638.]: # Logistic Regression Model
models.append('Logistic Regression')
log_reg = LogisticRegression(solver='liblinear', random_state=42)
start_time = time()
log_reg.fit(X_train, y_train)
train_time.append(time() - start_time)
y_train_preds = log_reg.predict(X_train)
y_test_preds = log_reg.predict(X_test)

Problem 9: Score the Model
What is the accuracy of your model?


In [639.]: train_score = log_reg.score(X_train, y_train)
train_accuracy.append(train_score)
test_accuracy.append(log_reg.score(X_test, y_test))
#accuracy scores
accuracy_score.append(metrics.accuracy_score(y_test, y_pred))
fpr, tpr, thresholds = metrics.roc_curve(y_test, y_pred)
AUC_score.append(metrics.auc(fpr, tpr))

In [640.]: from sklearn.metrics import confusion_matrix
confusion_matrix = confusion_matrix(y_test, y_test_preds)
print(confusion_matrix)

[[10803  178]
 [ 1140  354]]

In [641.]: disp = ConfusionMatrixDisplay(confusion_matrix=confusion_matrix)
disp.plot()

Out[641.]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at 0x16a38f6d0>



Problem 10: Model Comparisons
Now, we aim to compare the performance of the Logistic Regression model to our KNN algorithm, Decision Tree, and SVM models. Using the default settings for each of the models, fit and score each. Also, be sure to compare the fit time of each of the models. Present your findings in a DataFrame similar to that below:

Model  Train Time  Train Accuracy  Test Accuracy

RandomForestClassifier

In [64
```